

# MATH 1342

Section 5.3

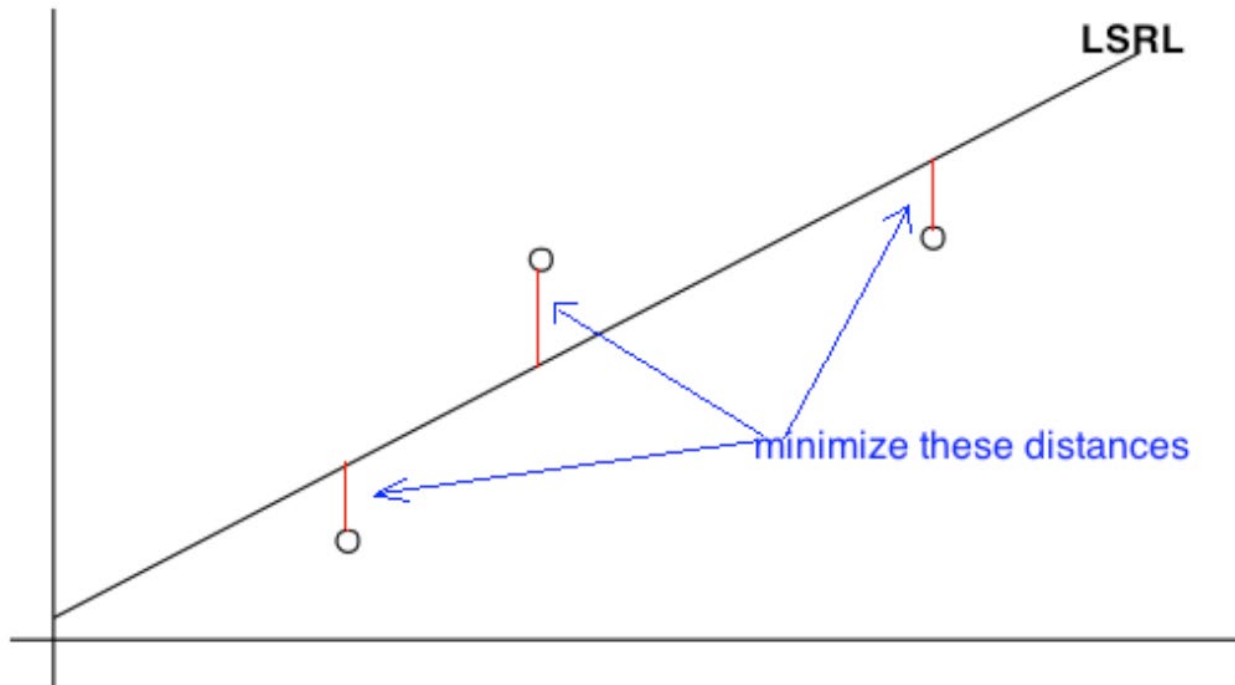
# Regression Lines

A **regression line** is a line that describes the relationship between the explanatory variable  $x$  and the response variable  $y$ .

Regression lines can be used to predict a value for  $y$  given a value of  $x$ .

# Least Squares Regression Lines (LSRL)

The **least squares regression line** (or LSRL) is a mathematical model used to represent data that has a linear relationship. We want a regression line that makes the vertical distances of the points in a scatter plot from the line as small as possible.



Note: To calculate this by hand, you are going to use optimization techniques from Calculus to minimize the distance between a point  $(x,y)$  from your scatter plot, and the line,  $y = mx + b$  by minimizing the distance formula:

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

# Calculating a Least Squares Regression Line

The least squares regression line formula is  $\hat{y} = a + bx$

The slope,  $b$  is calculated using  $b = r \left( \frac{s_y}{s_x} \right)$  and the  $y$ -intercept is  $a = \bar{y} - b\bar{x}$ .

To calculate the values of  $a$  and  $b$  for the regression line

with R-Studio, we use the command `lm(y ~ x)`

Example:

Using the Monopoly Problem, Calculate the Regression Line:

```
regline=lm(cost~spaces)
```

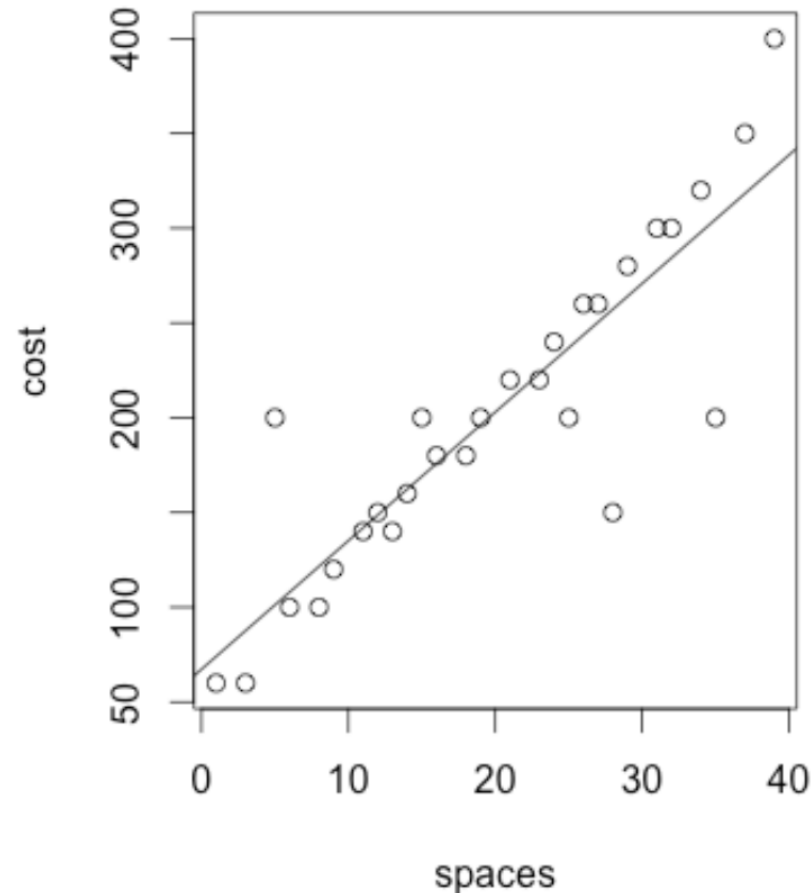
```
regline
```

<<This will give you the information about the linear equation>>

# Viewing the Scatterplot and the Regression Line

Note that I assigned a name to the `lm` command, this is not required unless you wish to use it again. We will use it again to plot the regression line on top of the scatterplot. The command is `abline`.

```
> abline(regline)
```



# Making Predictions:

The LSRL can be used to predict values of  $y$  given values of  $x$ .

Let's use our model to predict the cost of a property 50 spaces from GO

We need to be careful when predicting. When we are estimating  $y$  based on values of  $x$  that are much larger or much smaller than the rest of the data, this is called **extrapolation**.

# Interpreting the Slope

Notice that the formula for slope is  $b = r \left( \frac{s_y}{s_x} \right)$

this means that a change in one standard deviation in  $x$  corresponds to a change of  $r$  standard deviations in  $y$ . This means that on average, for each unit increase in  $x$ , there is an increase (or decrease if slope is negative) of  $|b|$  units in  $y$ .

Interpret the meaning of the slope for the Monopoly example



# Interpreting the Slope

Notice that the formula for slope is  $b = r \left( \frac{s_y}{s_x} \right)$

this means that a change in one standard deviation in  $x$  corresponds to a change of  $r$  standard deviations in  $y$ . This means that on average, for each unit increase in  $x$ , then is an increase (or decrease if slope is negative) of  $|b|$  units in  $y$ .

Interpret the meaning of the slope for the Monopoly example:

For every increase of 1 space from go, there is an increase of \$6.79 of cost.

# Coefficient of Determination

The square of the correlation ( $r$ ),  $r^2$  is called the **coefficient of determination**. It is the fraction of the variation in the values of  $y$  that is explained by the regression line and the explanatory variable.

When asked to interpret  $r^2$  we say, “approximately  $r^2 * 100\%$  of the variation in  $y$  is explained by the LSRL of  $y$  on  $x$ .”

This tells how accurate the measurement is based on the regression line.

# Facts about the coefficient of determination:

1. The coefficient of determination is obtained by squaring the value of the correlation coefficient.
2. The symbol used is  $r^2$
3. Note that  $0 \leq r^2 \leq 1$
4.  $r^2$  values close to 1 would imply that the model is explaining most of the variation in the dependent variable and *may be a very useful model*.
5.  $r^2$  values close to 0 would imply that the model is explaining little of the variation in the dependent variable and *may not be a useful model*.

Interpret  $r^2$  for the Monopoly problem

The following 9 observations compare the Quetelet index,  $x$  (a measure of body build) and dietary energy density,  $y$ .

$x$	221	228	223	211	231	215	224	233	268
$y$	.67	.86	.78	.54	.91	.44	.9	.94	.93

Compute the LSRL

Find the Correlation Coefficient

Find the coefficient of determination

Interpretation of the slope

Interpretation of the coefficient of determination